

Visualization of Reddit using Text Mining and Hierarchical Clustering

Per Blåwiik

Abstract—Based on the amount of available data from text based resources, it is relevant to explore potential usages for it. The purpose of this report was to present how to visualize and understand the website Reddit by utilizing common text mining methods, hierarchical clustering and force-directed graph algorithms. The fundamental idea was to use the most frequent words as basis for measuring similarities and subjectively interpret the content of subreddits. Using a force graph as a visual representation turned out to be intuitive, aesthetically pleasing and practical. Based on subjective analysis, the most frequent words deemed to be suitable both for clustering and understanding written text data.

Index Terms—Text mining, Hierarchical clustering, Force graph, Cosine similarity, Information visualization.

1 INTRODUCTION

Text mining, also referred as *text data mining*, can be briefly explained as a process where features or patterns are extracted from text based resources [1]. Examples of text based resources are the written content of websites such as online forums, customer reviews and articles. In the current age of social media, these types of resources are overflowing with available data and therefore it is relevant to explore the possible usages for it.

In order to understand abstract data and discover patterns, a graphical visualization tool is often necessary. Such a tool can include several steps: preparation of data by extracting important features, clustering data objects based on the features and visualize them in a graphical user interface by drawing various types of charts.

1.1 Problem

In this report, the implementation of a tool for visualizing the public online forum *Reddit* is explained. The content of reddit consists of posts made by the members of the website. A post consists of written text, links or images and can be up or down voted by other members. All posts belong to user-created topics called *subreddits*.

The purpose of the visualization tool is to help understand the contents of subreddits and find similarities between them. The following research questions were the basis for the project:

- Can the most frequent words in an arbitrary subreddit be used as an accurate description of the content of that subreddit?
- Can similarities between subreddits be identified when comparing their most frequent words? Can these words be used to represent a subreddit in order to cluster similar subreddits?
- Are there any correlations between the most frequent words and the given average scores of the texts where the words appeared?

The purpose of this report was to present methods that can help to answer the research questions with a visual representation of Reddit.

2 BACKGROUND AND RELATED WORK

A common application of text mining is the matching of search phrases (queries) against documents in a large database, based on relevance. For instance, if a database consists of medical journals, a doctor can form a query of keywords relevant to a particular syndrome. The query is then matched against the abstracts of all papers in the database and the system returns a list of documents ordered by a relevance ranking.

The query matching method is closely related to the second research question presented in Section 1.1: identify similarities between subreddits based on their most frequent words. For solving this problem, all posts in a subreddit must first be merged into a large text and then the words are counted. The top N most frequent words are then extracted into a vocabulary to represent that subreddit. Similarly to query matching, the subreddit vocabularies can be used to compute the relevance between them. Note that for query matching, a query is compared to a document collection. In the presented case, however, each document is compared to the document collection.

3 DATA

The data used for this report was derived from a dataset containing about 1.7 billion posts released by Reddit. All posts from the month of May 2015 (29.6 GB of data) is available online at *Kaggle*. Each post in the original dataset contains 22 features, however, only three of them were needed to answer the research questions: *Subreddit*, *Body* and *Score*. Subreddit is the name of the subreddit which the post belongs to; body contains the written text of the post; score is the number of up votes minus the number of down votes given to the post. An example of the raw format of the data is seen in Figure 1.

| Post | | |
|-----------|---|-------|
| Subreddit | Body | Score |
| soccer | Zlatan knows how to kick a ball! | 333 |
| AskReddit | I'm a guy and I had no idea this was a thing guys did. | -10 |
| music | Did well on my quiz today, am now eating ice cream. Good day. | 25 |

Fig. 1. The raw format of the data used for visualizing Reddit.

In the final implementation (before text preprocessing), 200 Megabytes of data were used. This corresponds to about 100 different subreddits, where each subreddit contains between 1000 - 20000 number of posts. These subreddits were subjectively chosen to cover a broad range of distinguished topics. The motivation for doing this was to make it easier to assess the success of clustering by relevance.

4 METHOD

In this section, the methods that were used for preprocessing, formatting, clustering and visualizing the data are presented.

4.1 Text Preprocessing

With the reduced dataset of 200 MB described in Section 3, several preprocessing steps were performed before the subreddits could be

The front-end system is the graphical user interface where users can interact with and explore the dataset. The force graph is the main tool for navigating this system. It offers the user to pan and zoom the view, and select subreddit nodes to understand their content. This system was developed as a *HTML* page using *JavaScript* and the library *D3*. A web browser based platform was used because it is supported by most computer systems. An additional advantage was that there were numerous of available examples for visualizing most types of data in various approaches using *D3*.

6 RESULTS

The final product was a visualization tool available for the most popular web browsers. The graphical user interface (front-end) is running on the client side and requires a connection to the server (back-end) to work. The result is an overview of 110 different subreddits laid out in a force graph where the nodes are coloured by a group id (generated from the hierarchical clustering). The most suitable number of clusters were subjectively found to be 12. The nodes spatial positioning is based on the direct similarities (cosine similarity) between the subreddits. A tolerance value of 0.09 was used to determine if two subreddits should be linked together. Each subreddit is approximated and represented by the top 50 most frequent words found in the posts belonging to the subreddit.

The main way of navigating the application is to interact with the force graph found in the center of the screen, see Figure 3. The user can pan and zoom the view of the force graph as well as selecting one or multiple nodes with the mouse. When selecting nodes, a list of the top 50 most frequent words is shown in a panel to the right side of the screen, see Figure 4.

The application responds instantaneously when selecting multiple subreddits since most of the computations were performed offline. All research questions presented in Section 1.1 could be answered by using the tool.

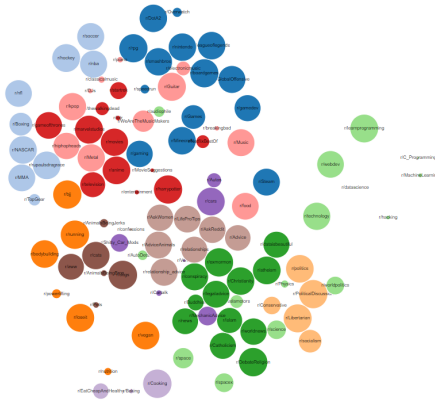


Fig. 3. The force graph visualizing all subreddits in the dataset. The node colours are based on hierarchical clustering, meaning that same coloured subreddits are related.

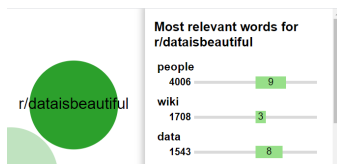


Fig. 4. A list of the top 50 most frequent words of the subreddit called *dataisbeautiful*. To the left below the words are the number of occurrences and to the right the average score.

7 EVALUATION

To determine if the visualization tool was intuitive and practical, an evaluation scheme was compiled. The aim of the evaluation was: to find out if the application was easy to understand and navigate, and, to learn if the visual representation suited the data.

7.1 Evaluation Scheme

Three evaluation methods were used: observation, think-out-loud and interview. By observing the actions of the participant when navigating the application, subconscious issues with the user interface can be revealed. When using the think-out-loud method, the interviewer can follow the thought process of the participant and discover subjective assumptions and frustration triggers. Finally, in the end, the participant can answer questions with their own words to highlight certain problems and offer suggestions for improving the system.

At the beginning of the session, the interviewer explained the dataset and the visual representation to the participant. Once the participant understood what they saw, they were given four tasks to complete. All tasks were aimed to force the participant to use the entire functionality of the application as well as to understand the visual representation of the data. While solving the tasks, the participant were asked to think-out-loud. The interviewer's role was to observe, take notes and intervene if the participant should get stuck. The participants had at all time access to the tasks in written form.

When the tasks were complete, the participants were given five questions to answer on a paper, asking: the purpose of the application, if anything was unclear or difficult, opinions on expected behaviour or responsiveness, and what they had learned from the data. At the end of the question form, they were also asked to make two statements on a scale from 1 to 4. An even scale was used to prevent neutral statements.

7.2 Evaluation Results

The evaluation was performed by five participants with no or little prior knowledge of Reddit. Note that this was just a coincident and not planned. The findings of the the evaluation was: the interface was easy to learn, the visual representation of the data was intuitive, the purpose of the tool was clear and the application was missing expected *quality-of-life* functionalities (e.g. search bar, single select multiple nodes).

8 CONCLUSIONS AND FUTURE WORK

All three of the research questions presented in Section 1.1 could, based on subjective analysis, be answered by using the final product: the most frequent words could describe the content of subreddits, similarities between subreddits was in many cases identified by using their most frequent words, the clustering of similar subreddits was successful and no correlation could be found between a words average score and its number of occurrences.

The force graph proved to be an intuitive visualization of the relevance between different subreddits (based on the nodes position and color) as well as offered a clear overview of the dataset. Panning and zooming the graph was fun and useful tools for exploring and understanding the dataset.

Some potential improvements for future work are to add a cluster bar to make it possible to select entire clusters, add a search bar for the force graph to select subreddits based on a search phrase and add a search bar for the list of the most frequent words.

REFERENCES

- [1] L. Eldén. *Matrix Methods in Data Mining and Pattern Recognition*, volume 2. SIAM, Philadelphia, PA, USA, 2019.
- [2] S. G. Kobourov. Spring embedders and force-directed graph drawing algorithms. January 2012.
- [3] L. Rokach and O. Maimon. Clustering methods. In *Data Mining and Knowledge Discovery Handbook*, pages 321–352, Boston, MA, USA, 2005. Springer.
- [4] J. H. Ward. Hierarchical grouping to optimize an objective function. *Journal of the American Statistical Association*, pages 236–244, 1963.